

Task 1

Create a database named 'custom'.

Create a table named temperature_data inside custom having below fields:

1. date (mm-dd-yyyy) format

2. zip code

3. temperature

The table will be loaded from comma-delimited file.

Load the dataset.txt (which is ',' delimited) in the table.

COMMAND:-

(i) create database custom;

(ii) create table temperature_data

```
(  
  dates STRING,  
  zip_code INT,  
  temperature SMALLINT  
)
```

row format delimited fields terminated by ',';

[note:- Later I use type casting over STRING data types of dates]

(iii) LOAD DATA LOCAL INPATH

'/home/acadgild/hive_folder/dataset_Session.txt' into table
 temperature_data;

(iv) select * from temperature_data;

[note:- here I cast the string data types of dates to date]

(v) select cast(to_date(from_unixtime(unix_timestamp(dates,
'dd-MM-yyyy')))) as date), zip_code, temperature from
 temperature_data;

EXPLANATION:-

(i) create the custom database here.

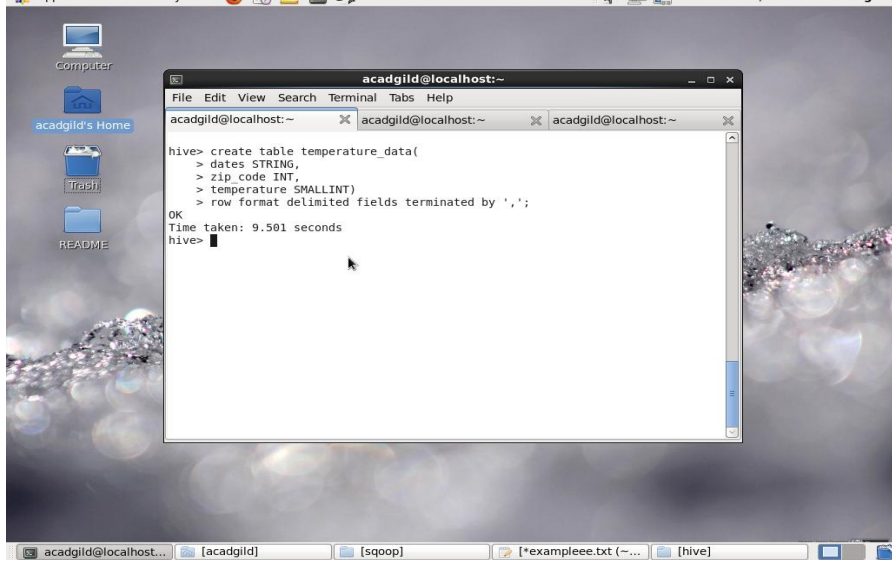
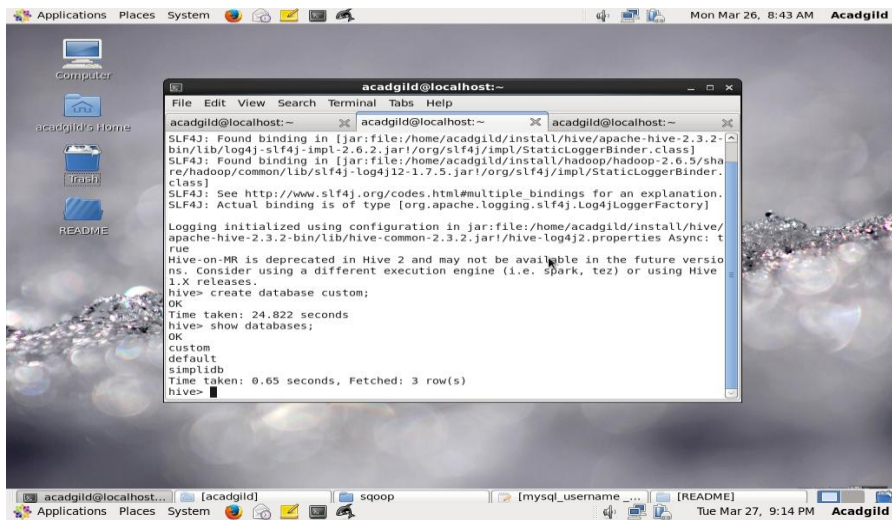
(ii) create the temperature_data here.

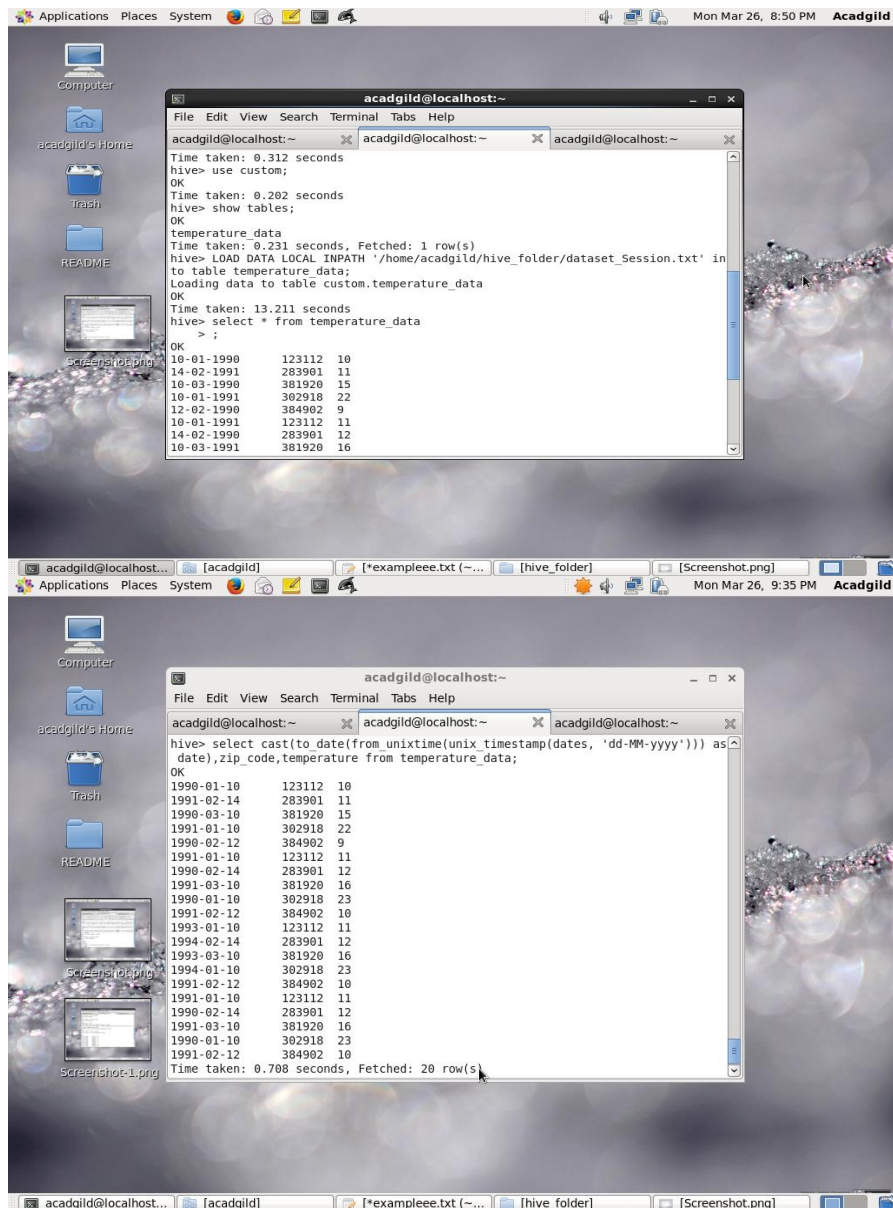
(iii) load the content of the dataset_Session.txt the temperature_data table.

(iv) check the temperature_data tables contents.

(v) after type casting check the contents of temperature_data again.

OUTPUT:-





Task 2

- Fetch date and temperature from temperature_data where zip code is greater than 300000 and less than 399999.
- Calculate maximum temperature corresponding to every year from temperature_data table.
- Calculate maximum temperature from temperature_data table corresponding to those years which have at least 2 entries in the table.
- Create a view on the top of last query, name it temperature_data_vw.
- Export contents from temperature_data_vw to a file in local file system, such that each file is '|' delimited.

COMMAND:-

(i) `select dates,temperature,zip_code from temperature_data where zip_code>300000 and zip_code<399999;`

(ii) `select YEAR(cast(to_date(from_unixtime(unix_timestamp(dates, 'dd-MM-yyyy')) as date)),MAX(temperature) from temperature_data group by YEAR(cast(to_date(from_unixtime(unix_timestamp(dates, 'dd-MM-yyyy')) as date)));`

(iii) `select YEAR(cast(to_date(from_unixtime(unix_timestamp(dates, 'dd-MM-yyyy')) as date)),MAX(temperature) from temperature_data group by YEAR(cast(to_date(from_unixtime(unix_timestamp(dates, 'dd-MM-yyyy')) as date)) having COUNT(YEAR(cast(to_date(from_unixtime(unix_timestamp(dates, 'dd-MM-yyyy')) as date)))>=2;`

(iv) `CREATE VIEW temperature_data_vw As select YEAR(cast(to_date(from_unixtime(unix_timestamp(dates, 'dd-MM-yyyy')) as date)),MAX(temperature) from temperature_data group by YEAR(cast(to_date(from_unixtime(unix_timestamp(dates, 'dd-MM-yyyy')) as date)) having COUNT(YEAR(cast(to_date(from_unixtime(unix_timestamp(dates, 'dd-MM-yyyy')) as date)))>=2;`

(v) `hive -S -e "USE custom; select * from temperature_data_vw" | sed 's/[\t]/|/g' > /home/acadgild/test.txt`

EXPLANATION:-

(i) Fetch date, temperature and zip_codes from temperature_data where zip code is greater than 300000 and less than 399999.

(ii) Calculate maximum temperature corresponding to every year from temperature_data table.

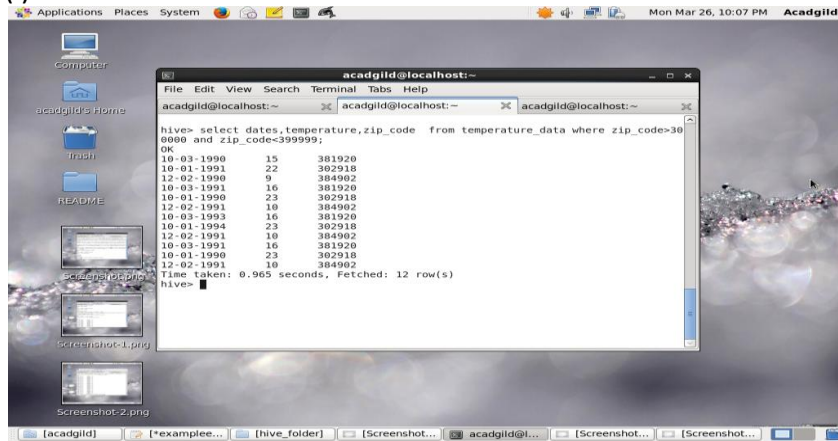
(iii) Calculate maximum temperature from temperature_data table corresponding to those years which have at least 2 entries in the table.

(iv) Create a view on the top of last query, name it temperature_data_vw.

(v) Export contents from temperature_data_vw to a file test.txt in local file system, such that each file is '|' delimited.

OUTPUT:-

(i)

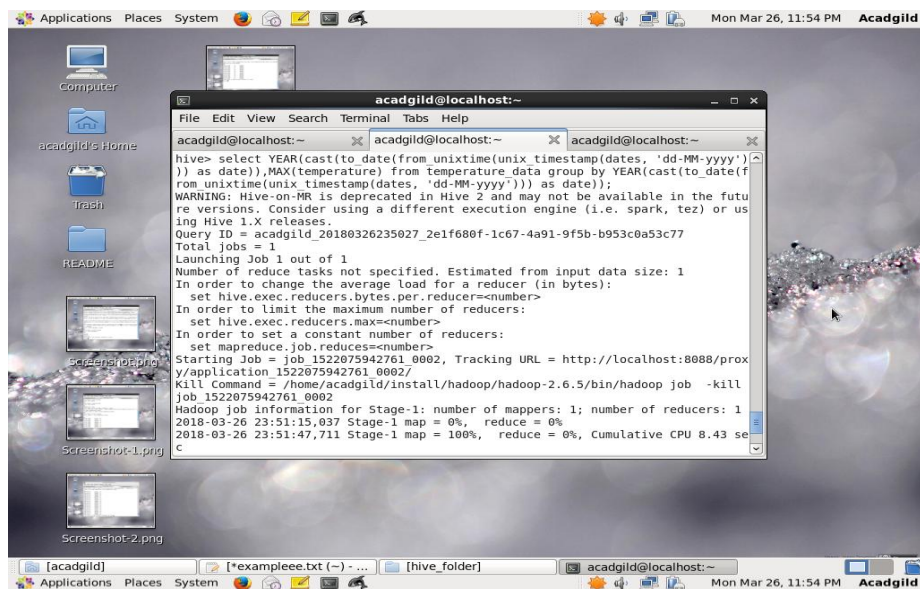


The screenshot shows a terminal window titled 'acadgild@localhost:~' with a menu bar (File, Edit, View, Search, Terminal, Tabs, Help). The terminal displays the output of a Hive query: 'select dates, temperature, zip_code from temperature_data where zip_code>30000 and zip_code<399999;'. The output is a table with three columns: dates, temperature, and zip_code. The data is as follows:

dates	temperature	zip_code
10-03-1990	15	381920
10-01-1991	22	382918
12-02-1990	9	384902
10-03-1991	16	381920
10-01-1990	23	382918
12-02-1991	10	384902
10-03-1993	16	381920
10-01-1994	23	382918
12-02-1991	10	384902
10-03-1991	16	381920
10-01-1990	23	382918
12-02-1991	10	384902

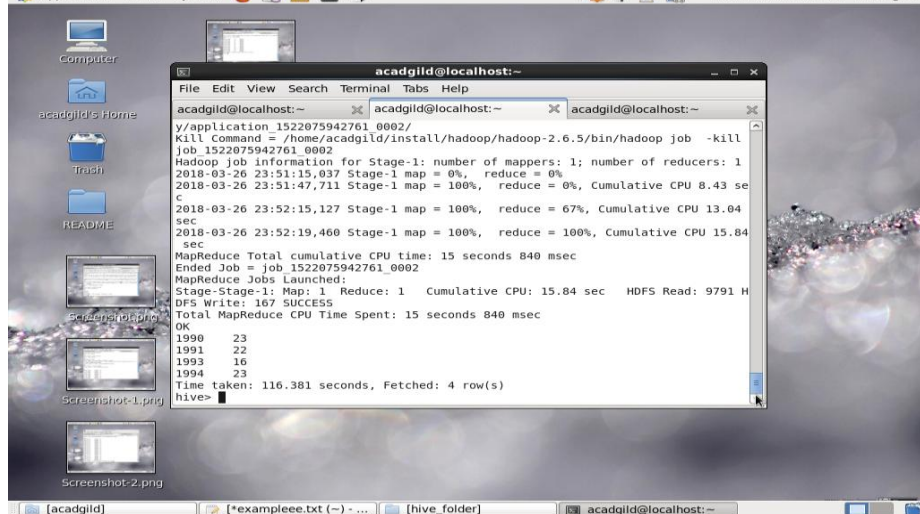
Time taken: 0.965 seconds, Fetched: 12 row(s)

(ii)



The screenshot shows a terminal window titled 'acadgild@localhost:~' with a menu bar (File, Edit, View, Search, Terminal, Tabs, Help). The terminal displays the output of a Hive query: 'select YEAR(cast(to date(from unixtime(unix timestamp(dates, 'dd-MM-yyyy')) as date)),MAX(temperature) from temperature_data group by YEAR(cast(to date(from unixtime(unix timestamp(dates, 'dd-MM-yyyy')) as date)))'. The output includes a warning about Hive-on-MR being deprecated, a query ID, and job execution details. The job is launched with 1 reducer. The output is as follows:

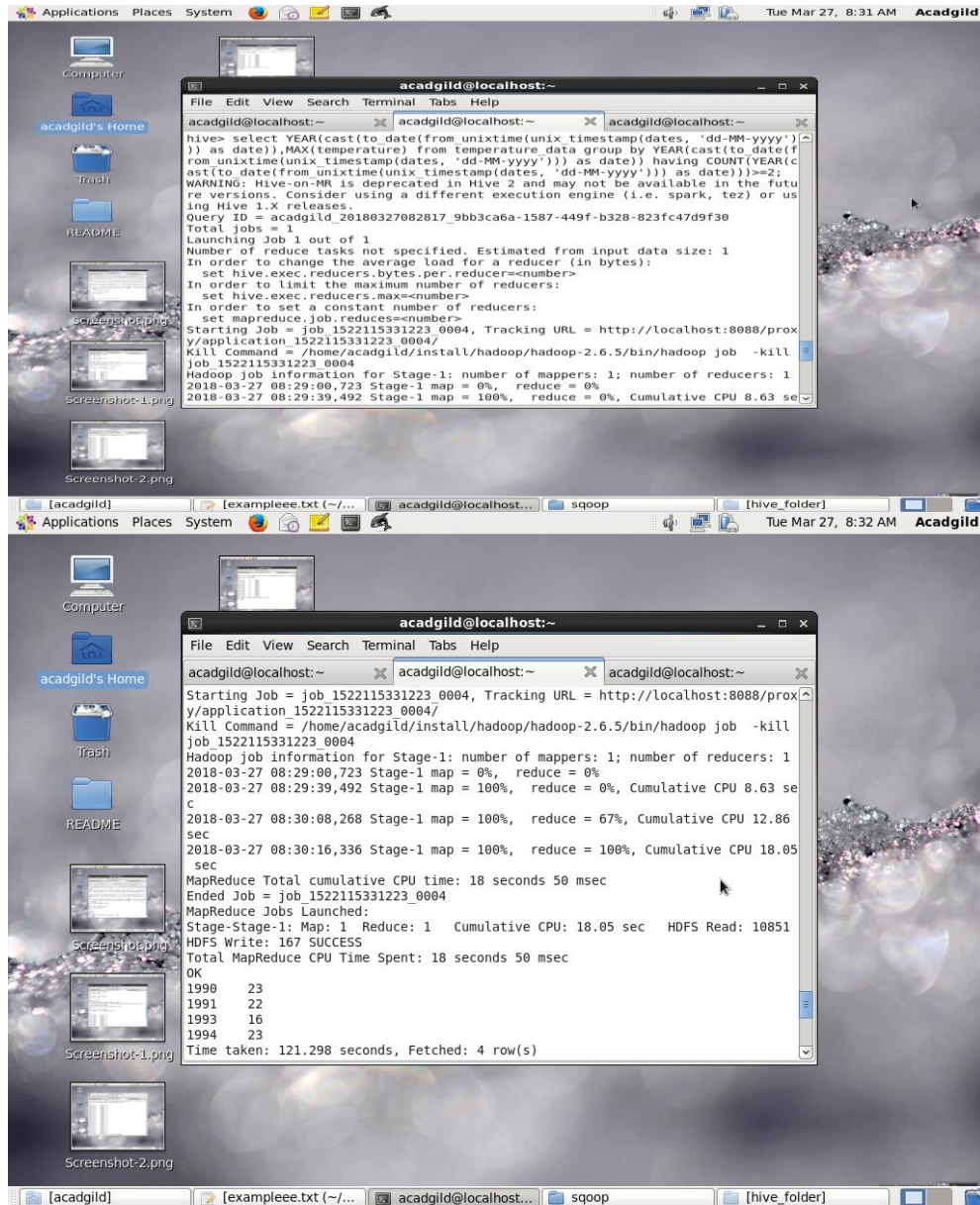
```
hives> select YEAR(cast(to date(from unixtime(unix timestamp(dates, 'dd-MM-yyyy')) as date)),MAX(temperature) from temperature_data group by YEAR(cast(to date(from unixtime(unix timestamp(dates, 'dd-MM-yyyy')) as date)))
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
Query ID = acadgild_20180326235027_2elf680f-1c67-4a91-9f5b-b953c0a53c77
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1522075942761_0002, Tracking URL = http://localhost:8088/proxy/application_1522075942761_0002/
Kill Command = /home/acadgild/install/hadoop/hadoop-2.6.5/bin/hadoop job -kill job_1522075942761_0002
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2018-03-26 23:51:15,037 Stage-1 map = 0%, reduce = 0%
2018-03-26 23:51:47,711 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 8.43 sec
```



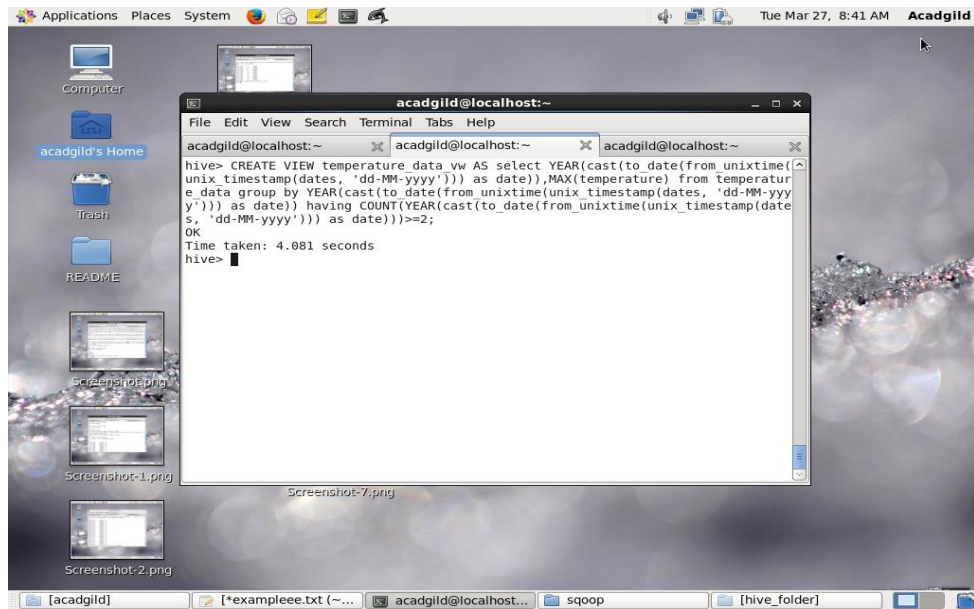
The screenshot shows a terminal window titled 'acadgild@localhost:~' with a menu bar (File, Edit, View, Search, Terminal, Tabs, Help). The terminal displays the output of a Hive query: 'select YEAR(cast(to date(from unixtime(unix timestamp(dates, 'dd-MM-yyyy')) as date)),MAX(temperature) from temperature_data group by YEAR(cast(to date(from unixtime(unix timestamp(dates, 'dd-MM-yyyy')) as date)))'. The output includes a warning about Hive-on-MR being deprecated, a query ID, and job execution details. The job is launched with 1 reducer. The output is as follows:

```
hives> select YEAR(cast(to date(from unixtime(unix timestamp(dates, 'dd-MM-yyyy')) as date)),MAX(temperature) from temperature_data group by YEAR(cast(to date(from unixtime(unix timestamp(dates, 'dd-MM-yyyy')) as date)))
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
Query ID = acadgild_20180326235027_2elf680f-1c67-4a91-9f5b-b953c0a53c77
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1522075942761_0002, Tracking URL = http://localhost:8088/proxy/application_1522075942761_0002/
Kill Command = /home/acadgild/install/hadoop/hadoop-2.6.5/bin/hadoop job -kill job_1522075942761_0002
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2018-03-26 23:51:15,037 Stage-1 map = 0%, reduce = 0%
2018-03-26 23:51:47,711 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 8.43 sec
2018-03-26 23:52:15,127 Stage-1 map = 100%, reduce = 67%, Cumulative CPU 13.04 sec
2018-03-26 23:52:19,460 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 15.84 sec
MapReduce Total cumulative CPU time: 15 seconds 840 msec
Ended Job = job_1522075942761_0002
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 15.84 sec HDFS Read: 9791 HDFS Write: 167 SUCCESS
Total MapReduce CPU Time Spent: 15 seconds 840 msec
OK
1990 23
1991 22
1993 16
1994 23
Time taken: 116.381 seconds, Fetched: 4 row(s)
hives>
```


(iii)



(iv)



(v)

